

A New QOS architecture to improve Scalability and Service in Heterogeneous Networks

¹*S.Senthilkumar*

Dept.Of CSE

Bharath university

²*Dr.M.Rajani*

Director Research

Bharath University

ABSTRACT

Our approach is simpler than prior network QOS schemes, which have required QOS support at every network node. In addition to reducing NOC area and energy consumption, the proposed topology-aware QOS architecture enables an elastic buffering (EB) optimization in parts of the network freed from QOS support. Elastic buffering further diminishes router buffer requirements by integrating storage into network links. We introduce a single-network EB architecture with lower cost compared to prior proposals. Our scheme combines elastic-buffered links and a small number of router-side buffers via a novel virtual channel allocation strategy. Our final NOC architecture is heterogeneous, employing QOS-enabled routers with conventional buffering in parts of the network, and light-weight elastic buffered nodes elsewhere. In a kilo-terminal NOC, this design enables a 29% improvement in power and a 45% improvement in area over a state-of-the-art QOS-enabled homogeneous network at the 15 nm technology node. In a modest-sized high-end chip, the proposed architecture reduces the NOC area to under 7% of the die and dissipates 23W of power when the network carries a 10% load factor averaged across the entire NOC. While the power consumption of the heterogeneous topology bests other approaches, low-energy CMPs and SOCs will be forced to better exploit physical locality to keep communication costs down. We further improve net-work area- and energy-efficiency through a novel flow control mechanism that enables a single-network, low-cost elastic buffer implementation. Together, these techniques yield a heterogeneous Kilo- NOC architecture that consumes 45% less area and 29% less power than a state-of-the-art QOS- enabled NOC without these features.

Key word: Design, Measurement, Performance

1. INTRODUCTION

Complexities of scaling single- with each technology generation, chips threaded performance have pushed processor containing over a thousand discrete designers in the direction of chip-level integration execution and storage resources will be likely in of multiple cores. Today's state-of-the-art general- the near future.

purpose chips integrate up to one hundred cores Chip-level multiprocessors (CMPs) require an while GPUs and other specialized processors may efficient communication infrastructure for contain hundreds of execution units. In addition operand, memory, coherence, and control transport, to the main processors, these chips often integrate motivating researchers to propose structured on-cache memories, specialized accelerators, memory chip networks as replacements to buses and ad- controllers, and other resources. Likewise, hoc wiring solutions of single- core chips. The modern systems-on-a-chip (SOCs) contains design of these networks-on- chip (NOCs) many cores, accelerators, memory channels, and typically re-quires satisfaction of multiple interfaces. As the degree of integration increases conflicting constraints, including minimizing

packet latency, reducing router area, and lowering communication energy overhead. In addition to basic packet transport, future NOCs will be expected to provide certain advanced services. In particular, quality-of-service (QOS) is emerging as a desirable feature due to the growing popularity of server consolidation, cloud computing, and real-time demands of SOCs. Despite recent advances aimed at improving the efficiency of individual NOC components such as buffers, crossbars, and flow control mechanisms, as well as features such as QOS, little attention has been paid to network scalability beyond several dozen terminals.

In this work, we focus on NOC scalability from the perspective of energy, area, performance, and quality-of-service. With respect to QOS, our interest is in mechanisms that provide hard guarantees, useful for enforcing Service Level Agreement (SLA) requirements in the cloud or real-time constraints in SOCs. Prior work showed that a direct low-diameter topology improves latency and energy efficiency in NOCs with dozens of nodes [16, 9]. While our analysis confirms this result, we identify critical scalability bottlenecks in these topologies once scaled to configurations with hundreds of network nodes. Chief among these is the buffer overhead associated with large credit round-trip times of long channels.

Large buffers adversely affect NOC area and energy efficiency. The addition of QOS support further increases storage overhead, virtual channel (VC) requirements, and arbitration complexity. For instance, a 256-node NOC with a low-diameter Multidrop Express Channel (MECS) topology and Preemptive Virtual Clock (PVC) QOS mechanism may require 750 VCs per router and over 12 MBs of buffering per chip.

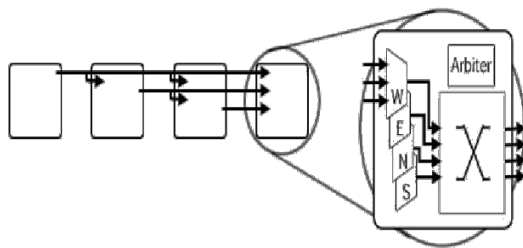


Figure 1: Multidrop Express Channel architecture.

In this paper, we propose a hybrid NOC architecture that offers low latency, small footprint, good energy efficiency, and SLA-strength QOS guarantees. The architecture is designed to scale to a large number of on-chip nodes and is evaluated in the context of a thousand terminal (Kilo-NOC) systems. To reduce the substantial QOS-related overheads, we address a key limitation of prior NOC QOS approaches which have required hardware support at every router node. Instead, our proposed topology-aware QOS architecture consolidates shared resources (e.g. memory controllers) within a portion of the network and only enforces QOS within sub-networks that contain these shared resources. The rest of the network, freed from the burden of hardware QOS support, enjoys diminished cost and complexity. Our approach relies on a richly-connected low-diameter topology to enable single-hop access to any QOS-protected sub-network, effectively eliminating intermediate nodes as sources of interference. To our knowledge, this work is the first to consider the interaction between topology and quality-of-service. Despite a significant reduction in QOS-related overheads, buffering remains an important contributor to our router area and energy footprint. We eliminate much of the expense by introducing a light-weight elastic buffer (EB) architecture that integrates storage directly into links, again using the topology to our advantage. To avoid deadlock in the resulting network, our approach leverages the multi-drop capability of a MECS interconnect to establish a dynamically allocated escape path for blocked packets into intermediate routers along the channel. In contrast, earlier EB schemes required multiple networks or many virtual channels for deadlock-free operation, incurring significant area and wire cost [21]. In a kilo-terminal network, the proposed single-network elastic buffer architecture requires only two virtual channels and reduces router storage requirements by 8x over a baseline MECS router without QOS support and by 12x compared to a QOS-enabled design.

Our results show that these techniques synergistically work to improve performance, area, and energy efficiency. In a kilo-terminal network in 15 nm technology, our final QOS-enabled NOC design reduces network area by 30% versus a modestly-provisioned MECS network with no QOS support and 45% compared to a MECS network with PVC, a prior NOC QOS architecture. Network energy efficiency improved by 29% and 40% over MECS without and with QOS support, respectively, on traffic with good locality. On random traffic, the energy savings diminish to 20% and 29% over the respective MECS baselines as

	Mesh	FBfly	MECS
Network diameter	$2k$	2	2
Bisection Channels / dimension	2	$K^2/2$	K
Crossbar (network ports)	4×4	$K \times k$	4×4
Arbitration	$\log(4)$	$\log(kv)$	$-\log(k-v)$

2. LITERATURE REVIEW

This section reviews key NOC concepts, draws on prior work to identify important Kilo-NOC technologies, and analyzes their scalability bottlenecks. We start with conventional NOC attributes – topology, flow control, and routing followed by quality-of-service technologies

2.1 Conventional NOC Attributes

2.1.1 Topology

Network topology determines the connectivity among nodes and is therefore a first-order determinant of network performance and energy-efficiency. To avoid the large hop counts associated with rings and meshes of early NOC designs, researchers have turned to richly-connected low-diameter networks that leverage the extensive on-chip wire budget. Such topologies reduce the number of costly router traversals at intermediate hops, thereby improving network latency and energy efficiency, and constitute a foundation for a Kilo-NOC.

One low-diameter NOC topology is the flattened butterfly (FBfly), which maps a richly – connected

wire energy dominates router energy consumption. Our NOC obtains both area and energy benefits without compromising either performance or QOS Guarantees.

In a notional 256MM² high-end chip, the proposed NOC consumes under 7% of the overall area and 23.5W of power at a sustained network load of 10%, a mod-est. fraction of the overall power budget.

Table 1: Scalability of NOC topologies. k : network radix, v : per-port VC count, C : a small integer

butterfly net-work to planar substrates by fully interconnecting nodes in each of the two dimensions via dedicated point-to-point channels. An alternative topology called Multidrop Express Channels (MECS) uses point-to- multipoint channels to also provide full intra-dimension connectivity but with fewer links. Each node in a MECS network has four out-put channels, one per cardinal direction. Lightweight drop interfaces allow packets to exit the channel into one of the routers spanned by the link. Figure 1 shows the high-level architecture of a MECS channel and router.

Scalability: Potential scalability bottlenecks in low-diameter networks are channels, input buffers, crossbar switches, and arbiters. The scaling trends for these structures are summarized in Table 1. The flattened butterfly requires $O(K^2)$ bisection channels per row/column, where K is the network radix, to support all-to-all intra-dimension connectivity. In contrast, the bisection channel count in MECS grows linearly with the radix.

Buffer capacities need to grow with network radix, assumed to scale with technology, to cover the round-trip credit latencies of long channel spans. Doubling the network radix doubles the number of input channels and the average buffer depth at an input port, yielding a quadratic increase in buffer capacity per node. This relationship holds for both flattened butterfly and MECS topologies and represents a true scalability obstacle.

Crossbar complexity is also quadratic in the number of input and output ports. This feature is

problematic in a flattened butterfly network, where port count grows in proportion to the network radix and causes a quadratic increase in switch area for every 2x increase in radix. In a MECS network, crossbar area stays nearly constant as the number of output ports is fixed at four and each switch input port is multiplexed among all network inputs from the same direction (see Figure 1). While switch complexity is not a concern in MECS, throughput can suffer because of the asymmetry in the number of input and output ports.

Finally, arbitration complexity grows logarithmically with port count. Designing a single-cycle arbiter for a high-radix router with a fast clock may be a challenge; however, arbitration can be pipelined over multiple cycles. While pipelined arbitration increases node delay, it is compensated for by the small hop count of low-diameter topologies. Hence, we do not consider arbitration a scalability bottleneck.

2.1.2 Flow Control

Flow control governs the flow of packets through the network by allocating channel bandwidth and buffer slots to packets. Conventional interconnects have traditionally employed packet-granularity bandwidth and storage allocation, exemplified by Virtual Cut-Through (VCT) flow control. In contrast, NOCs have relied on flit-level flow control, refining the allocation granularity to reduce the per-node storage requirements.

Scalability: In a Kilo-NOC with a low-diameter topology, long channel traversal times necessitate deep buffers to cover the round-trip credit latency. At the same time, wide channels reduce the number of flits per network packet. These two trends diminish the benefits of flit-level allocation since routers typically have enough buffer capacity for multiple packets. In contrast, packet-level flow control couples bandwidth and storage allocation, reducing the number of required arbiters, and amortizes the allocation delay over the length of a packet. Thus, in a Kilo-NOC, packet-level flow control is preferred to a flit-level architecture.

Elastic buffering: Recent research has explored the benefits of integrating storage elements, referred to as elastic buffers (EB), directly into network links. Kodi et al. proposed a scheme called iDEAL that augments conventional virtual-channel architecture with in-link storage, demonstrating savings in buffer area and power. An alternative proposal by Michelogiannakis et al. advocates a pure elastic-buffered architecture without any virtual channels. To prevent protocol deadlock in the resulting wormhole-routed NOC, the scheme requires a dedicated network for each packet class.

Scalability: To prevent protocol deadlock due to the serializing nature of buffered links, iDEAL must reserve a virtual channel at the destination router for each packet. As a result, its router buffer requirements in a low-diameter NOC grow quadratically with network radix as explained in Section 2.1.1, impeding scalability. A pure elastic-buffered architecture enjoys linear scaling in router storage requirements, but needs multiple networks for deadlock avoidance, incurring chip area and wiring expense.

2.1.3 Routing

A routing function determines the path of a packet from its source to the destination. Most networks use deterministic routing schemes, whose chief appeal is simplicity. In contrast, adaptive routing can boost throughput of a given topology at the cost of additional storage and/or allocation complexity.

Scalability: The scalability of a routing algorithm is a function of the path diversity attainable for a given set of channel resources. Compared to rings and meshes, direct low-diameter topologies typically offer greater path diversity through richer channel resources. Adaptive routing on such topologies has been shown to boost throughput [16, 9]; however, the gains come at the expense of energy efficiency due to the overhead of additional router traversals. While we do not consider routing a scalability bottleneck, reliability requirements may require additional

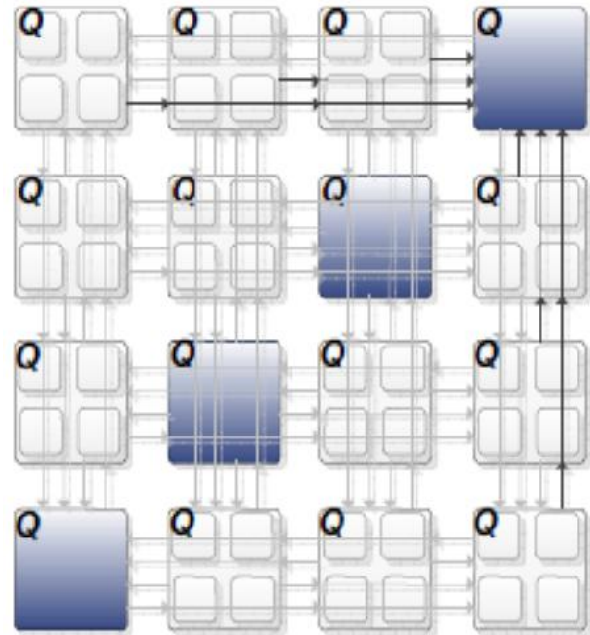
complexity not considered in this work.

2.2 Quality-of-Service

Cloud computing, server consolidation, and real-time applications demand on-chip QOS support for security, performance isolation, and guarantees. In many cases, a soft-ware layer will be unable to meet QOS requirements due to the fine-grained nature of chip-level resource sharing. Thus, we anticipate that hardware quality-of-service infrastructure will be a desirable feature in future CMPs. Unfortunately, existing network QOS schemes represent a weighty proposition that conflicts with the objectives of an area- and energy-scalable NOC. Current network QOS schemes require dedicated per-flow packet buffers at all network routers or source nodes, resulting in costly area and energy overheads. Recently proposed Preemptive Virtual Clock (PVC) architecture for NOC QOS relaxes the buffer requirements by using preemption to guarantee freedom from priority inversion. Under PVC, routers are provisioned with a minimum number of virtual channels (VCs) to cover the round-trip credit de-lay of a link. Without dedicated buffer resources for each flow, lower priority packets may block packets with higher dynamic priority. PVC detects such priority inversion situations and resolves them through preemption of lower- priority packets. Discarded packets require retransmission, signaled via a dedicated ACK network.

Scalability: While PVC significantly reduces QOS cost over prior work, in a low-diameter topology its VC requirements grow quadratic ally with network radix (analysis is similar to the one in Section 2.1.1), impeding scalability. VC requirements grow because multiple packets are not al-lowed to share a VC to prevent priority inversion within a FIFO buffer. Thus, longer links require more, but not deeper, VCs. Large VC populations adversely affect both storage requirements and arbitration complexity. In addition, PVC maintains per-flow state at each router whose storage requirements grow linearly with network size. Finally, preemption events in

PVC incur energy and latency overheads proportional to network diameter and preemption frequency. These considerations argue for an alternative net-work organization that provides QOS guarantees without compromising efficiency



(a) Baseline QOS-enabled CMP

(b) Topology-aware QOS approach Figure 1: 64-tile CMP with 4-way concentration and MECS topology. Light nodes: core+ cache tiles; shaded nodes: memory controllers; Q: QOS hardware. Dotted lines: domains in a topology-aware QOS architecture

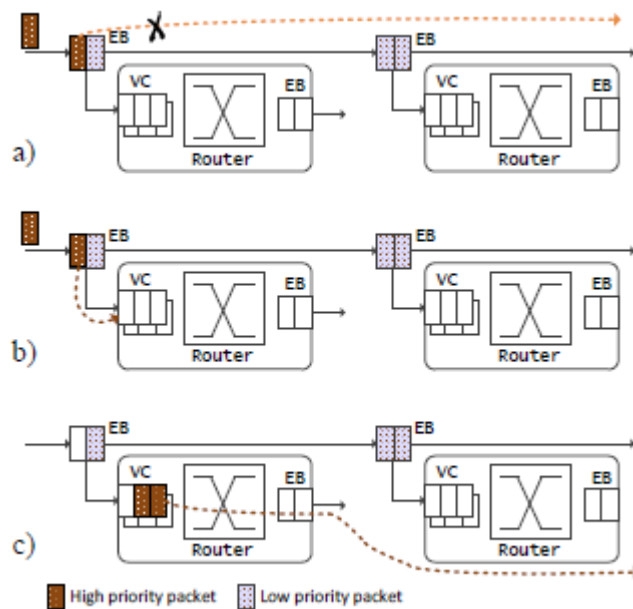
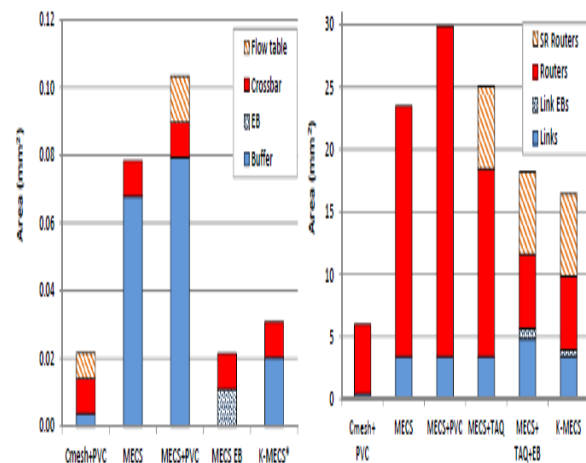


Figure2: Elastic buffer deadlock avoidance.

3. RESULTS & DISCUSSION

We first evaluate the different network organizations on area and energy-efficiency. Next, we compare the performance of elastic buffered networks to conventionally buffered designs. We then discuss QOS implications of various topologies. Finally, we examine performance stability and QOS on a collection of trace-driven workloads. Our area model accounts for four primary components of area overhead: input buffers, crossbar switch fabric, flow state tables, and router-side elastic buffers. Results are shown in Figure 3(a). The MECS+EB and K-MECS* bars corresponds to a router outside the shared region; all TAQ-enabled configurations use MECS+PVC routers inside the SR. We observe that elastic buffering is very effective in reducing router area in a MECS topology. Compared to a baseline MECS router with no QOS support, K-MECS* reduces router area by 61%. The advantage increases to 70% versus a PVC-enabled MECS router.



(a) Area of a single router (b) Total network area
Figure 3: Router and network area efficiency

4. CONCLUSION

In this paper, we proposed and evaluated architectures for kilo scale networks-on-chip (NOC) that address area, energy, and quality-of-service (QOS) challenges for large-scale on chip interconnects. We identify a low-diameter topology as a key Kilo-NOC technology for improving network performance and energy efficiency. While researchers have proposed low-diameter architectures for on-chip networks their scalability and QOS properties have not been studied. Our analysis reveals that large buffer requirements and QOS overheads stunt the ability of such topologies to support Kilo-NOC configurations in an area- and energy efficient fashion. We take a hybrid approach to network scalability. To reduce QOS overheads, we isolate shared resources in dedicated, QOS-equipped regions of the chip, enabling a reduction in router complexity in other parts of the die. The facilitating technology is a low-diameter topology, which affords single-hop interference free access to the QOS-protected regions from any node.

5. REFERENCES

- [1] J. D. Balfour and W. J. Dally. Design Tradeoffs for Tiled CMP On-chip Networks. In

International Conference on Supercomputing, pages 187–198, June 2006.

[2] C. Bienia, S. Kumar, J. P. Singh, and K. Li. The PARSEC Benchmark Suite: Characterization and Architectural Implications. In International Conference on Parallel Architectures and Compilation Techniques, pages 72–81, October 2008.

[3] N. L. Binkert, R. G. Dreslinski, L. R. Hsu, K. T. Lim, A. G. Saidi, and S. K. Reinhardt. The M5 Simulator: Modeling Networked Systems. IEEE Micro, 26(4):52–60, July/August 2006.

[4] W. J. Dally. Virtual-channel Flow Control. In International Symposium on Computer Architecture, pages 60–68, June 1990.

[5] W. J. Dally and B. Towles. Route Packets, Not Wires: On-chip Interconnection